

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2002-251197

(43)Date of publication of application : 06.09.2002

(51)Int.Cl.

G10L 15/00

G10L 15/06

G10L 15/10

G10L 15/16

G10L 17/00

H04N 5/91

(21)Application number : 2001-376561

(71)Applicant : NEC CORP

(22)Date of filing : 11.12.2001

(72)Inventor : IKOU KYOU

LIU XIN

(30)Priority

Priority number : 2000 254534

Priority date : 12.12.2000

Priority country : US

2001 011215

25.10.2001

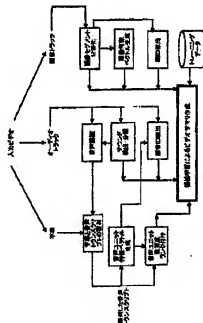
US

(54) AUDIOVISUAL SUMMARY CREATING METHOD

(57)Abstract:

PROBLEM TO BE SOLVED: To create an audio-centric, an image-centric, and an integrated audiovisual summaries of high quality by seamlessly integrating an image, audio, and text features extracted from input video.

SOLUTION: Integrated summarization is employed when strict synchronization of audio content and an image content is not required. A video programming which requires synchronization of the audio content and the image content is summarized by using an audio-centric approach or an image-centric approach. Both a machine learning-based approach and an alternative, heuristics-based approach are usable. Various probabilistic methods such as a naive Bayes method, a decision tree method, a neural network method, and a maximum entropy method are employed with the machine learning-based learning approach. To create an integrated audiovisual summary by using the alternative, heuristics-based approach, a maximum-bipartite-matching approach is employed.



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号
特開2002-251197
(P2002-251197A)

(43) 公開日 平成14年9月6日(2002.9.6)

(51) Int.Cl. ⁷	識別記号	F I	テラコード ⁸ (参考)
G 1 0 L 15/00		G 1 0 L 3/00	5 5 1 G 5 C 0 5 3
15/06			5 2 1 F 5 D 0 1 5
15/10			5 3 1 N
15/16			5 3 9
17/00			5 4 5 A
審査請求 有 請求項の数106 O L (全 19 頁) 最終頁に続く			
(21) 出願番号	特願2001-376561(P2001-376561)	(71) 出願人	000004237 日本電気株式会社 東京都港区芝五丁目7番1号
(22) 出願日	平成13年12月11日(2001.12.11)	(72) 発明者	イコウ キョウ アメリカ合衆国、ニュージャージー 08540 プリンストン、4 インディペン デンス ウエイ、エヌ・イー・シー・ユ ー・エス・エー インク内
(31) 優先権主張番号	6 0 / 2 5 4 5 3 4	(74) 代理人	100097157 弁理士 桂木 雄二
(32) 優先日	平成12年12月12日(2000.12.12)		
(33) 優先権主張国	米国 (U S)		
(31) 優先権主張番号	1 0 / 0 1 1 2 1 5		
(32) 優先日	平成13年10月25日(2001.10.25)		
(33) 優先権主張国	米国 (U S)		

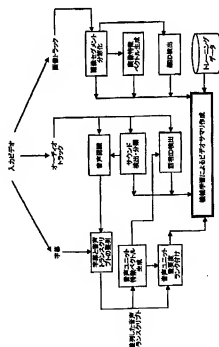
最終頁に続く

(54) 【発明の名称】 オーディオビジュアルサマリ作成方法

(57) 【要約】

【課題】 入力ビデオから抽出された画像、オーディオ、およびテキスト特徴をシームレスに統合することにより、オーディオ中心型、画像中心型、およびオーディオビジュアル統合型の高品質のサマリを作成する。

【解決手段】 オーディオと画像の内容の厳密な同期が要求されないときには、統合型サマリ作成が用いられる。オーディオ内容と画像内容の同期を要求するビデオ番組の場合、オーディオ中心型または画像中心型のいずれかの方法を用いてサマリが作成される。機械学習による方法と、代替法である発見的方法が使用可能である。ナイーブベイズ法、決定木法、ニューラルネットワーク法、および最大エントロピー法のようなさまざまな確率論的方法が、機械学習による方法で使用可能である。代替法である発見的方法を用いてオーディオビジュアル統合型サマリを作成するには、最大2部マッチング法が用いられる。



【特許請求の範囲】

【請求項1】 オーディオトラックおよび画像トラックを有するビデオ番組のオーディオ中心型オーディオビジュアルサマリを作成する方法において、前記オーディオビジュアルサマリの時間長 L_{sum} を選択するステップと、

前記オーディオトラックおよび画像トラックを検査するステップと、
前記オーディオビジュアルサマリの所望される内容に関連する1個以上の所定のオーディオ、画像、およびテキスト特性に基づいて、前記オーディオトラックから1個以上のオーディオセグメントを識別し、当該識別が、前記ビデオ番組内のオーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含まれるのに適している確率を与える、前もって生成された経験に基づく学習データに

依拠する機械学習法に従って実行される識別ステップと、
前記オーディオセグメントを前記オーディオビジュアルサマリに追加するステップと、
時間長 L_{sum} に達するまで、前記確率の降順に前記識別および追加を実行するステップと、

1個以上の識別されたオーディオセグメントに対応する1個以上の画像セグメントのみを、前記1個以上のオーディオセグメントと前記1個以上の画像セグメントの間の同期の程度が高くなるように、選択するステップとを有することを特徴とするオーディオビジュアルサマリ作成方法。

【請求項2】 前記識別するステップは、非音声サウンドを含むオーディオセグメントを検出するステップと、
内容に従って前記非音声サウンドを分類するステップと、

前記非音声サウンドのそれぞれについて、開始時刻コード、長さ、およびカテゴリを出力するステップとを有することを特徴とする請求項1記載の方法。

【請求項3】 前記オーディオセグメントが音声を含むとき、前記識別するステップは、前記オーディオセグメントに対する音声認識を実行して音声トランスクリプトを生成するステップと、
前記音声トランスクリプトのそれぞれについて、開始時刻コードおよび長さを出力するステップとを有することを特徴とする請求項2記載の方法。

【請求項4】 字幕が存在するとき、前記方法は、字幕と音声トランスクリプトを整理させるステップをさらに有することを特徴とする請求項3記載の方法。

【請求項5】 前記識別するステップは、前記字幕が存在する場合には前記整理に基づいて、また、前記字幕が存在しない場合には前記音声トランスクリプトに基づいて、音声ユニットを生成するステップ

と、
前記音声ユニットのそれぞれについて、特徴ベクトルを生成するステップとを有することを特徴とする請求項4記載の方法。

【請求項6】 前記音声ユニットのそれぞれについて、重要度ランクを計算するステップをさらに有することを特徴とする請求項5記載の方法。

【請求項7】 前記音声ユニットを受け取るステップと、

1以上の話者の識別を決定するステップとをさらに有することを特徴とする請求項6記載の方法。

【請求項8】 前記識別するステップは、前記画像トラックを個々の画像セグメントに分散化するステップを有することを特徴とする請求項1記載の方法。

【請求項9】 画像特徴を抽出するステップと、
前記画像セグメントのそれぞれについて、画像特徴ベクトルを形成するステップとをさらに有することを特徴とする請求項8記載の方法。

【請求項10】 前記画像セグメントのそれぞれについて、1個以上の顔の識別を決定するステップをさらに有することを特徴とする請求項9記載の方法。

【請求項11】 前記確率は、ナイーブベイズ法、決定木法、ニューラルネットワーク法、および最大エントロピー法からなる群から選択される方法に従って計算されることを特徴とする請求項1記載の方法。

【請求項12】 オーディオトラックおよび画像トラックを有するビデオ番組の画像中心型オーディオビジュアルサマリを作成する方法において、

前記オーディオビジュアルサマリの時間長 L_{sum} を選択するステップと、
前記ビデオ番組の前記画像トラックおよびオーディオトラックを検査するステップと、

前記オーディオビジュアルサマリの所望される内容に関連する1個以上の所定の画像、オーディオ、およびテキスト特性に基づいて、前記画像トラックから1個以上の画像セグメントを識別し、当該識別が、前記ビデオ番組内の前記画像セグメントのそれぞれについて、与えられた画像セグメントが前記オーディオビジュアルサマリに

含まれるのに適している確率を与える、前もって生成された経験に基づく学習データに依拠する機械学習法に従って実行される識別ステップと、
前記1個以上の画像セグメントを前記オーディオビジュアルサマリに追加するステップと、
時間長 L_{sum} に達するまで、前記確率の降順に前記識別および追加を実行するステップと、

1個以上の識別された画像セグメントに対応する1個以上のオーディオセグメントのみを、前記1個以上の画像セグメントと前記1個以上のオーディオセグメントの間の同期の程度が高くなるように、選択するステップとを有することを特徴とするオーディオビジュアルサマリ作

成方法。

【請求項13】 前記識別するステップは、前記画像トラックを個々の画像セグメントに分割化するステップを有することを特徴とする請求項12記載の方法。

【請求項14】 画像特徴を抽出するステップと、前記画像セグメントのそれぞれについて、画像特徴ベクトルを形成するステップとをさらに有することを特徴とする請求項13記載の方法。

【請求項15】 前記画像セグメントのそれぞれについて、1個以上の顔の識別を決定するステップをさらに有することを特徴とする請求項10記載の方法。

【請求項16】 前記オーディオビジュアルサマリ内の前記画像セグメントのそれぞれについて、最小再生時間 L_{min} を選択するステップをさらに有することを特徴とする請求項12記載の方法。

【請求項17】 比較的多数のオーディオセグメントおよび画像セグメントが前記オーディオビジュアルサマリに提供されて、幅指向のオーディオビジュアルサマリを提供するように、 L_{min} は L_{max} に比べて十分に小さいことを特徴とする請求項16記載の方法。

【請求項18】 比較的小数のオーディオセグメントおよび画像セグメントが前記オーディオビジュアルサマリに提供されて、深さ指向のオーディオビジュアルサマリを提供するように、 L_{min} は L_{max} に比べて十分に大きいことを特徴とする請求項16記載の方法。

【請求項19】 前記識別するステップは、非音声サウンドを含むオーディオセグメントを抽出するステップと、

内容に従って前記非音声サウンドを分類するステップと、前記非音声サウンドのそれぞれについて、開始時刻コード、長さ、およびカテゴリを出力するステップとを有することを特徴とする請求項12記載の方法。

【請求項20】 前記オーディオセグメントが音声を含むとき、前記識別するステップは、前記オーディオセグメントに対する音声認識を実行して音声トランスクリプトを生成するステップと、前記音声トランスクリプトのそれぞれについて、開始時刻コードおよび長さを出力するステップとを有することを特徴とする請求項19記載の方法。

【請求項21】 字幕が存在するとき、前記方法は、字幕と音声トランスクリプトを整理させるステップをさらに有することを特徴とする請求項20記載の方法。

【請求項22】 前記識別するステップは、前記字幕が存在する場合には前記整理に基づいて、また、前記字幕が存在しない場合には前記音声トランスクリプトに基づいて、音声ユニットを生成するステップと、前記音声ユニットのそれぞれについて、特徴ベクトルを生成するステップとを有することを特徴とする請求項20

1記載の方法。

【請求項23】 前記音声ユニットのそれぞれについて、重要度ランクを計算するステップをさらに有することを特徴とする請求項22記載の方法。

【請求項24】 前記音声ユニットを受け取るステップと、1以上の話者の識別を決定するステップとをさらに有することを特徴とする請求項23記載の方法。

【請求項25】 前記確率は、ナイーブベイズ法、決定木法、ニューラルネットワーク法、および最大エントロピー法からなる群から選択される方法に従って計算されることを特徴とする請求項12記載の方法。

【請求項26】 オーディオトラックおよびビデオトラックを有するビデオ番組の統合オーディオビジュアルサマリを作成する方法において、前記オーディオビジュアルサマリ内の時間長 L_{max} を選択するステップと、

オーディオビジュアルサマリに含まれるべき前記画像セグメントのそれぞれについて、最小再生時間 L_{min} を選択するステップと、

前記オーディオビジュアルサマリ内の長さ L_{max} に達するまで1個以上の所望されるオーディオセグメントを選択し、当該選択が、前記ビデオ番組内の前記オーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含まれるのに適している確率を与える、前もって生成された経験に基づく学習データに依拠する機械学習法に従って実行されることによりオーディオサマリを作成するステップと、

前記画像セグメントのそれぞれについて、前記機械学習法に従って、与えられた画像セグメントが前記オーディオビジュアルサマリに含まれるのに適している確率を計算するステップと、

選択された前記オーディオセグメントのそれぞれについて、対応する画像セグメントに関して、得られるオーディオセグメントと画像セグメントのペアが、あらかじめ規定された整列要件を満たすかどうかを検査するステップと、

得られるオーディオセグメントと画像セグメントのペアが、あらかじめ規定された整列要件を満たす場合、ペアのオーディオセグメントと画像セグメントをそれぞれの最初から前記最小再生時間 L_{min} だけ整列させて第1の整列点を規定するステップと、

前記検査および識別を繰り返して前記整列点をすべて識別するステップと、

前記オーディオビジュアルサマリ内の全長を、前記オーディオビジュアルサマリ内の最初から開始し第1の整列点で終了するか、

ある整列点における画像セグメントの最後から開始し次の整列点で終了するか、

最後の整列点における画像セグメントの最後から開始し前記オーディオビジュアルサマリの中で終了するかのいずれかの期間をそれぞれ有する複数のパーティションに分割するステップと、
前記パーティションの期間にそれぞれについて、以下のステップ、すなわち、
該パーティションの期間に入る画像セグメントの集合を識別するステップと、
前記パーティションに挿入されることが可能な画像セグメントの個数を決定するステップと、
挿入されるべきと識別された画像セグメントの長さを決定するステップと、
与えられた画像セグメントが前記オーディオビジュアルサマリに含まれるのに適している前記確率の降順に、識別された画像セグメントを前記個数だけ選択するステップと、
選択された画像セグメントのそれぞれについて、それぞれの最初から前記時間長だけのセクションを収集し、すべての収集されたセクションを時間の降順に前記パーティションに追加するステップとに従って、さらに画像セグメントを追加するステップとを有することを特徴とするオーディオビジュアルサマリ作成方法。

【請求項27】 前記識別するステップは、非音声サウンドを含むオーディオセグメントを検出するステップと、
内容に従って前記非音声サウンドを分類するステップと、
前記非音声サウンドのそれぞれについて、開始時刻コード、長さ、およびカテゴリを出力するステップとを有することを特徴とする請求項26記載の方法。

【請求項28】 前記オーディオセグメントが音声を含むとき、前記識別するステップは、
前記オーディオセグメントに対する音声認識を実行して音声トランスクリプトを生成するステップと、
前記音声トランスクリプトのそれぞれについて、開始時刻コードおよび長さを出力するステップとを有することを特徴とする請求項27記載の方法。

【請求項29】 字幕が存在するとき、前記方法は、字幕と音声トランスクリプトを整列させるステップをさらに有することを特徴とする請求項28記載の方法。

【請求項30】 前記字幕が存在する場合には前記整列に基づいて、また、前記字幕が存在しない場合には前記音声トランスクリプトに基づいて、音声ユニットを生成するステップと、
前記音声ユニットのそれぞれについて、特徴ベクトルを生成するステップとをさらに有することを特徴とする請求項29記載の方法。

【請求項31】 前記音声ユニットのそれぞれについて、重要度ランクを計算するステップをさらに有することを特徴とする請求項30記載の方法。

【請求項32】 前記音声ユニットを受け取るステップと、
1以上の話者の識別を決定するステップとをさらに有することを特徴とする請求項31記載の方法。

【請求項33】 比較的多数の画像セグメントが前記オーディオビジュアルサマリに提供されて、幅指向のオーディオビジュアルサマリを提供するように、 L_{ms} は L_{ms} に比べて十分に小さいことを特徴とする請求項26記載の方法。

10 【請求項34】 比較的小数の画像セグメントが前記オーディオビジュアルサマリに提供されて、深さ指向のオーディオビジュアルサマリを提供するように、 L_{ms} は L_{ms} に比べて十分に大きいことを特徴とする請求項26記載の方法。

【請求項35】 前記与えられたオーディオセグメントが前記オーディオビジュアルサマリに含まれるのに適している確率は、ナイーブベイズ法、決定木法、ニューラルネットワーク法、および最大エントロピー法からなる群から選択される方法に従って計算されることを特徴とする請求項26記載の方法。

20 【請求項36】 前記与えられた画像セグメントが前記オーディオビジュアルサマリに含まれるのに適している確率は、ナイーブベイズ法、決定木法、ニューラルネットワーク法、および最大エントロピー法からなる群から選択される方法に従って計算されることを特徴とする請求項26記載の方法。

【請求項37】 前記識別するステップは、前記画像トラックを個々の画像セグメントに分散化するステップを有することを特徴とする請求項26記載の方法。

30 【請求項38】 画像特徴を抽出するステップと、
前記画像セグメントのそれぞれについて、画像特徴ベクトルを形成するステップとをさらに有することを特徴とする請求項37記載の方法。

【請求項39】 前記画像セグメントのそれぞれについて、1個以上の顔の識別を決定するステップをさらに有することを特徴とする請求項38記載の方法。

【請求項40】 オーディオトラックおよび画像トラックを有するビデオ番組のオーディオ中心型オーディオビジュアルサマリを作成する方法において、

40 前記オーディオビジュアルサマリの時間長 L_{ms} を選択するステップと、
前記オーディオトラックおよび画像トラックを検査するステップと、

前記オーディオビジュアルサマリの所望される内容に関連する1個以上の所定のオーディオ、画像、およびテキスト特性に基づいて、前記オーディオトラックから1個以上のオーディオセグメントを識別し、当該識別が、前記ビデオ番組内の前記オーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含まれるのに適しているか

どうかを決定するランク付けを与える、所定の発見的ルールの集合に従って実行される識別ステップと、前記オーディオセグメントを前記オーディオビジュアルサマリに追加するステップと、

時間長 L_{max} に達するまで、前記オーディオセグメントのランク付けの降順に前記識別および追加を実行するステップと、

1個以上の識別されたオーディオセグメントに対応する1個以上の画像セグメントのみを、前記1個以上のオーディオセグメントと前記1個以上の画像セグメントの間の同期の程度が高くなるように、選択するステップとを有することを特徴とするオーディオビジュアルサマリ作成方法。

【請求項41】 前記識別するステップは、非音声サウンドを含むオーディオセグメントを検出するステップと、内容に従って前記非音声サウンドを分類するステップと、前記非音声サウンドのそれぞれについて、開始時刻コード、長さ、およびカテゴリを出力するステップとを有することを特徴とする請求項40記載の方法。

【請求項42】 前記オーディオセグメントが音声を含むとき、前記識別するステップは、前記オーディオセグメントに対する音声認識を実行して音声トランスクリプトを生成するステップと、前記音声トランスクリプトのそれぞれについて、開始時刻コードおよび長さを出力するステップとを有することを特徴とする請求項41記載の方法。

【請求項43】 字幕が存在するとき、前記方法は、字幕と音声トランスクリプトを整合させるステップをさらに有することを特徴とする請求項42記載の方法。

【請求項44】 前記字幕が存在しない場合には前記整合に基づいて、また、前記字幕が存在しない場合には前記音声トランスクリプトに基づいて、音声ユニットを生成するステップと、前記音声ユニットのそれぞれについて、特徴ベクトルを生成するステップとをさらに有することを特徴とする請求項43記載の方法。

【請求項45】 前記音声ユニットを受け取るステップと、

1以上の話者の識別を決定するステップとをさらに有することを特徴とする請求項44記載の方法。

【請求項46】 前記識別するステップは、前記画像トラックを個々の画像セグメントに分割化するステップを有することを特徴とする請求項40記載の方法。

【請求項47】 画像特徴を抽出するステップと、前記画像セグメントのそれぞれについて、画像特徴ベクトルを形成するステップとをさらに有することを特徴とする請求項46記載の方法。

【請求項48】 前記画像セグメントのそれぞれについ

て、1個以上の顔の識別を決定するステップをさらに有することを特徴とする請求項47記載の方法。

【請求項49】 前記音声ユニットのそれぞれについて前記ランク付けを計算するステップをさらに有することを特徴とする請求項40記載の方法。

【請求項50】 オーディオトラックおよび画像トラックを有するビデオ番組の画像中心型オーディオビジュアルサマリを作成する方法において、前記サマリの時間長 L_{max} を選択するステップと、前記画像トラックおよびオーディオトラックを検査するステップと、

前記オーディオビジュアルサマリの所望される内容に関連する1個以上の所定の画像、オーディオ、およびテキスト特性に基づいて、前記画像トラックから1個以上の画像セグメントを識別し、当該識別が、前記ビデオ番組内の前記画像セグメントのそれぞれについて、与えられた画像セグメントが前記オーディオビジュアルサマリに含められるのに適しているかどうかを決定するランク付けを与える、所定の発見的ルールの集合に従って実行される識別ステップと、

前記1個以上の画像セグメントを前記オーディオビジュアルサマリに追加するステップと、時間長 L_{max} に達するまで、前記ランク付けの降順に前記識別および追加を実行するステップと、1個以上の識別された画像セグメントに対応する1個以上のオーディオセグメントのみを、前記1個以上の画像セグメントと前記1個以上のオーディオセグメントの間の同期の程度が高くなるように、選択するステップとを有することを特徴とするオーディオビジュアルサマリ作成方法。

【請求項51】 前記識別するステップは、所定の視覚的類似性および動的特性に基づいて、前記ビデオ番組の画像セグメントをクラス化するステップを有することを特徴とする請求項50記載の方法。

【請求項52】 前記識別するステップは、前記画像トラックを個々の画像セグメントに分割化するステップを有することを特徴とする請求項51記載の方法。

【請求項53】 画像特徴を抽出するステップと、前記フレームクラスターのそれぞれについて、画像特徴ベクトルを形成するステップとをさらに有することを特徴とする請求項52記載の方法。

【請求項54】 前記フレームクラスターのそれぞれについて、1個以上の顔の識別を決定するステップをさらに有することを特徴とする請求項53記載の方法。

【請求項55】 前記識別するステップは、非音声サウンドを含むオーディオセグメントを検出するステップと、内容に従って前記非音声サウンドを分類するステップと、前記非音声サウンドのそれぞれについて、開始時刻コ

ド、長さ、およびカテゴリを出力するステップとを有することを特徴とする請求項50記載の方法。

【請求項56】 前記オーディオセグメントが音声を含むとき、前記識別するステップは、前記オーディオセグメントに対する音声認識を実行して音声トランスクリプトを生成するステップと、前記音声トランスクリプトのそれぞれについて、開始時刻コードおよび長さを出力するステップとを有することを特徴とする請求項55記載の方法。

【請求項57】 字幕が存在するとき、前記方法は、字幕と音声トランスクリプトを整理させるステップをさらに有することを特徴とする請求項56記載の方法。

【請求項58】 前記字幕が存在する場合には前記整理に基づいて、また、前記字幕が存在しない場合には前記音声トランスクリプトに基づいて、音声ユニットを生成するステップと、

前記音声ユニットのそれぞれについて、特徴ベクトルを生成するステップとをさらに有することを特徴とする請求項57記載の方法。

【請求項59】 前記音声ユニットのそれぞれについて、重要度ランクを計算するステップをさらに有することを特徴とする請求項58記載の方法。

【請求項60】 前記音声ユニットを受け取るステップと、

1以上の話者の識別を決定するステップとをさらに有することを特徴とする請求項59記載の方法。

【請求項61】 前記オーディオビジュアルサマリ内の前記画像セグメントのそれぞれについて、最小再生時間 L_{min} を選択するステップをさらに有することを特徴とする請求項50記載の方法。

【請求項62】 比較的多数のオーディオセグメントおよび画像セグメントが前記オーディオビジュアルサマリに提供されて、幅指向のオーディオビジュアルサマリを提供するように、 L_{min} は L_{min} に比べて十分に小さいことを特徴とする請求項61記載の方法。

【請求項63】 比較的小数のオーディオセグメントおよび画像セグメントが前記オーディオビジュアルサマリに提供されて、深さ指向のオーディオビジュアルサマリを提供するように、 L_{min} は L_{min} に比べて十分に大きいことを特徴とする請求項61記載の方法。

【請求項64】 オーディオトラックおよびビデオトラックを有するビデオ番組の統合オーディオビジュアルサマリを作成する方法において、

前記オーディオビジュアルサマリの長さ L_{sum} を選択するステップと、

オーディオビジュアルサマリに含まれるべき複数の画像セグメントのそれぞれについて、最小再生時間 L_{min} を選択するステップと、

前記ビデオ番組内の前記オーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記

オーディオビジュアルサマリに含まれるのに適しているかどうかを決定するランク付けを与える、所定の発見のルールの集合に従って、1個以上の所望されるオーディオセグメントを選択することによって、オーディオサマリを作成するステップと、

前記オーディオビジュアルサマリの長さに達するまで、前記オーディオセグメントのランク付けの降順に、前記選択を実行するステップと、

各フレームクラスが少なくとも1つの前記画像セグメントを含み、与えられたフレームクラス内のすべての画像セグメントが互いに視覚的に類似しているように、前記画像セグメントの視覚的類似性および動的特性に基づいて、前記ビデオ番組の前記画像セグメントを複数のフレームクラスへとグループ分けするステップと、

選択された前記オーディオセグメントのそれぞれについて、対応する画像セグメントに関して、得られるオーディオセグメントと画像セグメントのペアが、あらかじめ規定された整理要件を満たすかどうかを検査するステップと、

得られるオーディオセグメントと画像セグメントのペアが、あらかじめ規定された整理要件を満たす場合、ペアのオーディオセグメントと画像セグメントをそれぞれの最初から前記最小再生時間 L_{min} だけ整列させて第1の整列点を規定するステップと、前記検査および識別を繰り返して前記整列点をすべて識別するステップと、

前記オーディオビジュアルサマリの全長を、前記オーディオビジュアルサマリの最初から開始し第1の整列点で終了するか、

ある整列点における画像セグメントの最後から開始し次の整列点で終了するか、

最後の整列点における画像セグメントの最後から開始し前記オーディオビジュアルサマリの最後で終了するか、いずれかの期間をそれぞれ有する複数のパーティションに分割するステップと、

各時間スロットが前記最小再生時間 L_{min} に等しい長さを有するように、前記パーティションのそれぞれを複数の時間スロットに分割するステップと、

前記フレームクラスと前記時間スロットの間の最適マッチングに従って、以下のこと、すなわち、

各フレームクラスをただ1つの時間スロットに割り当てること、および、オーディオビジュアルサマリ内のすべての画像セグメントの時間順序を維持することに基づいて、前記パーティションのそれぞれの前記時間スロットを満たすように前記フレームクラスを割り当ててステップとを有することを特徴とするオーディオビジュアルサマリ作成方法。

【請求項65】 前記最適マッチングは、最大2部マッチング法によって計算されることを特徴とする請求項64記載の方法。

【請求項66】 フレームクラスタより多くの時間スロットがある場合、複数の画像セグメントを含むフレームクラスタを識別し、前記オーディオビジュアルサマリ内の前記画像セグメントの時間順序を維持しながら、すべての前記時間スロットが満たされるまで、前記識別されたフレームクラスタからの画像セグメントを時間スロットに割り当てることを特徴とする請求項65記載の方法。

【請求項67】 前記時間順序が維持されていることを確認するために前記オーディオビジュアルサマリを検査するステップと、前記時間順序が維持されていない場合、前記時間順序が維持されるように、各パーティションに追加された前記画像セグメントを並べ替えるステップとをさらに有することを特徴とする請求項66記載の方法。

【請求項68】 前記識別するステップは、非音声サウンドを含むオーディオセグメントを検出するステップと、内容に従って前記非音声サウンドを分類するステップと、前記非音声サウンドのそれぞれについて、開始時刻コード、長さ、およびカテゴリを出力するステップとを有することを特徴とする請求項64記載の方法。

【請求項69】 前記オーディオセグメントが音声を含むとき、前記識別するステップは、前記オーディオセグメントに対する音声認識を実行して音声トランスクリプトを生成するステップと、前記音声トランスクリプトのそれぞれについて、開始時刻コードおよび長さを出力するステップとを有することを特徴とする請求項68記載の方法。

【請求項70】 字幕が存在するとき、前記方法は、字幕と音声トランスクリプトを整理させるステップをさらに有することを特徴とする請求項69記載の方法。

【請求項71】 前記字幕が存在する場合には前記整理に基づいて、また、前記字幕が存在しない場合には前記音声トランスクリプトに基づいて、音声ユニットを生成するステップと、前記音声ユニットのそれぞれについて、特徴ベクトルを生成するステップとをさらに有することを特徴とする請求項70記載の方法。

【請求項72】 前記音声ユニットのそれぞれについて、重要度ランクを計算するステップをさらに有することを特徴とする請求項71記載の方法。

【請求項73】 前記音声ユニットを受け取るステップと、

1以上の話者の識別を決定するステップとをさらに有することを特徴とする請求項72記載の方法。

【請求項74】 比較的小数の画像セグメントが前記オーディオビジュアルサマリに提供されて、幅指向のオーディオビジュアルサマリを提供するように、 L_{min} は L

に比べて十分に小さいことを特徴とする請求項64記載の方法。

【請求項75】 比較的小数の画像セグメントが前記オーディオビジュアルサマリに提供されて、深さ指向のオーディオビジュアルサマリを提供するように、 L_{min} は L_{max} に比べて十分に大きいことを特徴とする請求項64記載の方法。

【請求項76】 前記識別するステップは、前記画像トラックを個々の画像セグメントに分散化するステップを有することを特徴とする請求項64記載の方法。

【請求項77】 画像特徴を抽出するステップと、前記フレームクラスタのそれぞれについて、画像特徴ベクトルを形成するステップとをさらに有することを特徴とする請求項76記載の方法。

【請求項78】 前記画像セグメントのそれぞれについて、1個以上の顔の識別を決定するステップをさらに有することを特徴とする請求項77記載の方法。

【請求項79】 オーディオビジュアルコンテンツからなるビデオ番組のビデオサマリを作成する装置におい

て、前記オーディオビジュアルコンテンツのオーディオトラックおよび画像トラックを検査する検査手段と、

前記ビデオサマリの所望されるコンテンツに関連する所定のオーディオ、画像、およびテキスト特性のうちの少なくとも1つに基づき、前記オーディオトラックから1個以上のオーディオセグメントを、前記画像トラックから1個以上の画像セグメントを、前記ビデオサマリに含められるのに適しているかどうかを決定する順位を与える所定の基準に従って、それぞれ識別する手段と、

前記順位に従って、前記1個以上のオーディオセグメントおよび前記1個以上の画像セグメントをそれぞれ時間軸上に配置して前記ビデオサマリを生成する手段と、を有することを特徴とするビデオサマリ作成装置。

【請求項80】 前記識別する手段は、前記ビデオサマリの所望される内容に関連する1個以上の所定のオーディオ、画像、およびテキスト特性に基づいて、前記オーディオトラックから1個以上のオーディオセグメントを識別する際に、当該識別を、前記ビデオ番組内のオーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含められるのに適している確率を与える、前もって生成された経験に基づく学習データに依拠する機械学習法に従って実行する、ことを特徴とする請求項79記載のビデオサマリ作成装置。

【請求項81】 前記識別する手段は、前記オーディオセグメントをカテゴリ化することを特徴とする請求項80記載のビデオサマリ作成装置。

【請求項82】 前記オーディオセグメントは、音声と非音声とにカテゴリ化されることを特徴とする請求項81記載のビデオサマリ作成装置。

【請求項83】 前記識別する手段は、非音声サウンドを含むオーディオセグメントを検出し、内容に従って前記非音声サウンドを分類し、前記非音声サウンドのそれぞれについて、オーディオ情報を出力することを特徴とする請求項82記載のビデオサマリ作成装置。

【請求項84】 前記オーディオ情報は、開始時刻コード、長さ、およびカテゴリであることを特徴とする請求項83記載のビデオサマリ作成装置。

【請求項85】 前記識別する手段は、前記オーディオビジュアルサマリに所望される内容に関連する1個以上の所定の画像、オーディオ、およびテキスト特性に基づいて、前記画像トラックから1個以上の画像セグメントを識別し、当該識別が、前記ビデオ番組内の前記画像セグメントのそれぞれについて、与えられた画像セグメントが前記オーディオビジュアルサマリに含められるのに適している確率を与える、前もって生成された経験に基づく学習データに依拠する機械学習法に従って実行する、ことを特徴とする請求項79記載のビデオサマリ作成装置。

【請求項86】 前記識別する手段は、1個以上の所望されるオーディオセグメントを選択し、当該選択が、前記ビデオ番組内の前記オーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含められるのに適している確率を与える、前もって生成された経験に基づく学習データに依拠する機械学習法に従って実行する、ことを特徴とする請求項79記載のビデオサマリ作成装置。

【請求項87】 前記識別する手段は、前記オーディオビジュアルサマリに所望される内容に関連する1個以上の所定のオーディオ、画像、およびテキスト特性に基づいて、前記オーディオトラックから1個以上のオーディオセグメントを識別し、当該識別が、前記ビデオ番組内の前記オーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含められるのに適しているかどうかを決定するランク付けを与える、所定の発見のルール集合に従って実行する、ことを特徴とする請求項79記載のビデオサマリ作成装置。

【請求項88】 前記識別する手段は、前記オーディオセグメントをカテゴリ化することを特徴とする請求項87記載のビデオサマリ作成装置。

【請求項89】 前記オーディオセグメントは、音声と非音声とにカテゴリ化されることを特徴とする請求項88記載のビデオサマリ作成装置。

【請求項90】 前記識別する手段は、非音声サウンドを含むオーディオセグメントを検出し、内容に従って前記非音声サウンドを分類し、前記非音声サウンドのそれぞれについて、オーディオ情報を出力す

ることを特徴とする請求項89記載のビデオサマリ作成装置。

【請求項91】 前記オーディオ情報は、開始時刻コード、長さ、およびカテゴリであることを特徴とする請求項90記載のビデオサマリ作成装置。

【請求項92】 前記識別する手段は、前記オーディオビジュアルサマリに所望される内容に関連する1個以上の所定の画像、オーディオ、およびテキスト特性に基づいて、前記画像トラックから1個以上の画像セグメントを識別し、当該識別が、前記ビデオ番組内の前記画像セグメントのそれぞれについて、与えられた画像セグメントが前記オーディオビジュアルサマリに含められるのに適しているかどうかを決定するランク付けを与える、所定の発見のルール集合に従って実行する、ことを特徴とする請求項79記載のビデオサマリ作成装置。

【請求項93】 オーディオビジュアルコンテンツからなるビデオ番組のビデオサマリを作成する方法において、

- 20 前記オーディオビジュアルコンテンツのオーディオトラックおよび画像トラックを検査し、前記ビデオサマリに所望されるコンテンツに関連する所定のオーディオ、画像、およびテキスト特性のうち少なくとも1つに基づき、前記オーディオトラックから1個以上のオーディオセグメントを、前記画像トラックから1個以上の画像セグメントを、前記ビデオサマリに含められるのに適しているかどうかを決定する順位を与える所定の基準に従って、それぞれ識別し、前記順位に従って、前記1個以上のオーディオセグメントおよび前記1個以上の画像セグメントをそれぞれ時間軸上に配置して前記ビデオサマリを生成する、ステップを有することを特徴とするビデオサマリ作成方法。

【請求項94】 前記識別するステップは、前記ビデオサマリに所望される内容に関連する1個以上の所定のオーディオ、画像、およびテキスト特性に基づいて、前記オーディオトラックから1個以上のオーディオセグメントを識別する際に、当該識別を、前記ビデオ番組内のオーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含められるのに適している確率を与える、前もって生成された経験に基づく学習データに依拠する機械学習法に従って実行する、ことを特徴とする請求項93記載のビデオサマリ作成方法。

【請求項95】 前記識別するステップは、前記オーディオセグメントをカテゴリ化することを特徴とする請求項94記載のビデオサマリ作成方法。

【請求項96】 前記オーディオセグメントは、音声と非音声とにカテゴリ化されることを特徴とする請求項95記載のビデオサマリ作成方法。

【請求項97】 前記識別するステップは、非音声サウンドを含むオーディオセグメントを検出し、内容に従って前記非音声サウンドを分類し、前記非音声サウンドのそれぞれについて、オーディオ情報を出力することを特徴とする請求項97記載のビデオサマリ作成方法。

【請求項98】 前記オーディオ情報は、開始時刻コード、長さ、およびカテゴリであることを特徴とする請求項97記載のビデオサマリ作成方法。

【請求項99】 前記識別するステップは、前記オーディオビジュアルサマリの所望される内容に関連する1個以上の所定の画像、オーディオ、およびテキスト特性に基づいて、前記画像トラックから1個以上の画像セグメントを識別し、当該識別が、前記ビデオ番組内の前記画像セグメントのそれぞれについて、与えられた画像セグメントが前記オーディオビジュアルサマリに含められるのに通じている確率を与える、前もって生成された経験に基づく学習データに依拠する機械学習法に従って実行する、ことを特徴とする請求項93記載のビデオサマリ作成方法。

【請求項100】 前記識別するステップは、1個以上の所望されるオーディオセグメントを選択し、当該選択が、前記ビデオ番組内の前記オーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含められるのに通じている確率を与える、前もって生成された経験に基づく学習データに依拠する機械学習法に従って実行する、ことを特徴とする請求項93記載のビデオサマリ作成方法。

【請求項101】 前記識別するステップは、前記オーディオビジュアルサマリの所望される内容に関連する1個以上の所定のオーディオ、画像、およびテキスト特性に基づいて、前記オーディオトラックから1個以上のオーディオセグメントを識別し、当該識別が、前記ビデオ番組内の前記オーディオセグメントのそれぞれについて、与えられたオーディオセグメントが前記オーディオビジュアルサマリに含められるのに通じているかどうかを決定するランク付けを与える、所定の発見のルールの集合に従って実行する、ことを特徴とする請求項93記載のビデオサマリ作成方法。

【請求項102】 前記識別するステップは、前記オーディオセグメントをカテゴリ化することを特徴とする請求項101記載のビデオサマリ作成方法。

【請求項103】 前記オーディオセグメントは、音声と非音声とにカテゴリ化されることを特徴とする請求項102記載のビデオサマリ作成方法。

【請求項104】 前記識別するステップは、非音声サウンドを含むオーディオセグメントを検出し、内容に従って前記非音声サウンドを分類し、前記非音声サウンドのそれぞれについて、オーディオ情報を出力す

ることを特徴とする請求項103記載のビデオサマリ作成方法。

【請求項105】 前記オーディオ情報は、開始時刻コード、長さ、およびカテゴリであることを特徴とする請求項104記載のビデオサマリ作成方法。

【請求項106】 前記識別するステップは、前記オーディオビジュアルサマリの所望される内容に関連する1個以上の所定の画像、オーディオ、およびテキスト特性に基づいて、前記画像トラックから1個以上の画像セグメントを識別し、当該識別が、前記ビデオ番組内の前記画像セグメントのそれぞれについて、与えられた画像セグメントが前記オーディオビジュアルサマリに含められるのに通じているかどうかを決定するランク付けを与える、所定の発見のルールの集合に従って実行する、ことを特徴とする請求項93記載のビデオサマリ作成方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、一般に、ビデオサマリ作成技術に関し、特に、入力ビデオから抽出した画像、オーディオ、およびテキスト特徴をシームレスに統合することによりビデオサマリを作成する方法およびシステムに関する。

【0002】

【従来の技術】長い論説や学術論文などのテキスト文書の多くには要約がある。要約の助けにより、読者は、文書全体を詳細に分析せずに、その文書の内容が関心のあるものかどうかをすばやく確かめることができる。テキスト文書の場合もそうであるが、ビデオ番組の内容および性質は一目では捉えられないことが多い。同様にして一般的に全体の内容を示すために、長いビデオ番組の要約すなわちサマリを提供することが一般に望まれる。

【0003】最近、ワールドワイドウェブ(WWWあるいはウェブ)の爆発的な成長により、オンラインテキストおよびマルチメディアデータコレクションの数が急激に増大している。オンラインマルチメディアコンテンツの増大というこの傾向が続くと、ユーザが大量のデータから最も関連性のある情報をすばやく識別することを支援する自動データサマリ作成技術はますます重要になる。

【0004】この状況において、ビデオサマリ作成が、困難な課題を提示する。その作業が困難であるのは、ビデオ番組の画像トラックおよびオーディオトラックの両方のサマリ作成をまず必要とするからである。2つのサマリを自然なやりかたで有効に統合することが、もう1つの課題となる。

【0005】一般に、ほとんどの種類のビデオサマリ作成は、オーディオ中心型サマリ作成(audio-centric summarization)、画像中心型サマリ作成(image-centric summarization)、およびオーディオビジュアル統合型サマリ

リ作成(integrated audio-visual summarization)という3つのカテゴリに分類することができる。ビデオ番組のうちには、例えばニュース放送、ドキュメンタリー、ビデオセミナーのように、対応するオーディオトラックと画像トラックの間に強い相関のないタイプのものがある。このようなビデオカテゴリについては、オーディオと画像をゆるく整列(整合)させながら、オーディオと画像の両方の内容のカバレッジを最大にするオーディオビジュアル統合型サマリ作成アプローチを使用するのが適当である。他方、映画、ドラマ、トークショーなどのような他のタイプのビデオ番組は、オーディオトラックと画像トラックの間に強い相関を有することがある。この種のビデオ番組については、オーディオ提示とビデオ画像の間の同期が重要である。このような状況では、オーディオ中心型または画像中心型のいずれかのサマリ作成方法を使用するのが適当である。

【0006】

【発明が解決しようとする課題】従来のシステムは、このようなさまざまなタイプのビデオ番組に対する有効で効率的なサマリ作成という課題に対し、包括的な解決法を提供していない。現在使用されている多くのビデオサマリ作成システム・方法は、あるタイプのビデオ内容を発見的に重要であるとし、これらのあらかじめ指定した内容を入力ビデオから抽出することによってサマリを作成している。その結果、これらのビデオサマリ作成システム・方法は、非常に領域特異的(領域固有)かつアプリケーション特異的であり、ユーザの個々の需要に基づいてサマリを作成することや、さまざまな種類のビデオ番組を処理することができない。

【0007】

【課題を解決するための手段】本発明は、機械学習フレームワークに基づいてビデオサマリ作成のシステムおよび方法を提供することによって、従来のビデオサマリ作成技術の前記およびその他の欠点を克服する。また、本発明はさらに、機械学習フレームワークによって要求されるトレーニングデータを得ることが困難な状況に対処するためのシステムおよび方法も提供する。これらのシステムおよび方法は、入力ビデオから抽出される画像、オーディオ、およびテキスト特徴をシステムに統合することによって、高品質のオーディオおよび画像のサマリを作成することができる。

【0008】オーディオトラックと画像トラックの間の強い同期に依存しないビデオ番組の具体例として、オーディオセグメントが最近の地震による犠牲者の数に関する情報を提示しているテレビニュース番組を考える。対応する画像セグメントは、現場のレポートの接写であったり、崩壊した建物の現場で作業する救助隊の接写であったり、地震の震央を示す地域地図の接写であったりする。このような場合、オーディオ内容は、必ずしも、対応する画像内容に言及している必要がないことが多

い。前述のように、このようなビデオ番組のその他の例には、ドキュメンタリー、セミナーなどがある。

【0009】本発明の一実施例によれば、厳密な同期が要求されないときには、ビデオ番組のサマリを作成するために、オーディオビジュアル統合型サマリ作成技術が用いられる。このようなビデオ番組のサマリを作成する際には、オーディオおよび画像のサマリを別個に作成することが好ましい。その後、2つのサマリが、ゆるく整列して統合される。このアプローチでは、オーディオ内容と画像内容の両方のカバレッジを、サマリにおいて最大化することが可能である。

【0010】逆に、オーディオ内容と画像内容の間の強い同期を要求するビデオ番組は、一般に、特定の瞬間におけるオーディオトラックがその瞬間に提示される画像と直接関係しており、その逆も同様であるということによって特徴づけられる。このようなビデオ番組のサマリを作成する際には、オーディオと画像の間の同期が重要である。したがって、同期はオーディオ中心型または画像中心型のいずれかであることが好ましい。

20 【0011】一実施例によれば、オーディオ中心型サマリ作成技術は、ビデオ番組に関連するオーディオ内容の重要な側面を確認する。必要な程度の同期を達成するため、画像サマリは、オーディオサマリを構成するオーディオセグメントに対応するビデオフレームを選択することによってのみ、生成される。画像中心サマリ作成技術は、まず、ビデオ番組の重要な画像セグメントを識別することによって画像トラックのサマリを作成する。その後、これらの重要なあるいは代表的な画像セグメントに対応するオーディオセグメントを、全体のビデオサマリに含める。

30 【0012】サマリを作成するプロセスは、画像、オーディオ信号、音声トランスクリプト、および字幕(クロードキャプション)テキストからの手がかりおよび特徴を利用することによって容易化される。画像特徴、音声トランスクリプト、および字幕テキストは、オーディオサマリ作成を改善するために、対応するオーディオ特徴と組み合わせられ、一方、オーディオ特徴、音声トランスクリプト、および字幕テキストは、よりよい画像サマリ作成を容易にするために、関連する画像特徴と組み合わせられる。

40 【0013】オーディオ中心型、画像中心型、あるいはオーディオビジュアル統合型のサマリ作成を実現するため、以下では2つの実施例について説明する。1つの技術によれば、与えられたアプリケーションに対していずれのサマリ作成技術が好ましいかに応じて、あらかじめサマリ作成の選択(プレファレンス)を例示することが可能なトレーニングデータを用いて、機械学習が、ビデオ番組のオーディオあるいは画像トラックに適用される。この技術では、システムは、既知のアルゴリズム方式のうちの任意のものを用いて、サンプルビデオサマリ

に示される挙動を模倣し、このサンプルから、および、サンプルの固有のインプリメンテーションから、学習を行うことが可能である。必要な命令をシステムに提供するために、トレーニングデータが直ちに入手可能でない場合や容易に適用可能でない場合には、以下で説明するもう1つの実施例が、本発明の代替方法として、適用可能である。

【0014】本発明の上記および関連するその他の利点は、添付図面を参照して、以下の好ましい実施例の詳細な説明を検討すれば、さらに明らかとなる。

【0015】

【発明の実施の形態】図面を参照すると、図1は、機械学習によるビデオサマリ作成システム・方法に関する、本発明の一実施例の動作を示す流れ図である。図1を参照して、以下では、使用される数学的モデルのタイプ、オーディオおよびビジュアルサマリ作成に用いられる特徴、ならびに、オーディオおよびビジュアルサマリを整理させる方法について説明する。

【0016】【機械学習フレームワーク】通常のビデオ番組は、オーディオトラックおよび画像トラックの両方を含み、これらはいずれも長く連続することがある。このようなビデオ番組のサマリを作成するには、そのビデオを構成するオーディオトラックおよび画像トラックの両方を、意味かつ管理可能な操作ユニットに分類化しなければならない。例えば、意味あるオーディオ操作ユニットとしては、1個の単語、1個の句、1個の文、あるいはその他のコヒーレントな音響プロファイルを含むオーディオセグメントの発声がある。同様に、可能な画像操作ユニットの例には、単一のカメラショット、一連の連続するカメラショット、ある判断基準によってグループ分けされた画像フレームのクラスなどがあ

る。【0017】このような状況において、あるベクトルすなわち特徴セットXで、オーディオまたは画像操作ユニットを表すことが可能である。さらに、Xは、いくつかの特徴xを含む。特徴xは、オーディオまたは画像操作ユニットに関連する画像特徴、オーディオ特徴、テキスト特徴（例えば、音声トランスクリプトや字幕からの重要なキーワード）とすることが可能である。n個の特徴xが特定のベクトルすなわち特徴セットXに存在する場合、 $X = [x_1, x_2, \dots, x_n]$ である。サマリ作成作業は、与えられた特徴セットXに対して、確率 $P(y|X)$ を計算する二分類問題に変換される。ここでyは2進（バイナリ）変数であり、その値1および0は、Xがサマリに含まれるか否かのそれぞれの状態を表す。この確率 $P(y|X)$ は、ルール（規則）によるアプローチを用いて決定することも可能であり、あるいは、機械学習法を用いて評価することも可能である。後者の場合、トレーニングデータが機械学習システムに提供され、システムは、提供されたトレーニングデータに従って、確率 $P(y|X)$ を予測するモデルを学習すること

になる。

【0018】確率 $P(y|X)$ を評価するために、ナイーブベイズ法、決定木法、ニューラルネットワーク法、最大エントロピー法（これらには限定されない）などのような、既知のさまざまな機械学習技術のうちの任意のものを使用可能である。このような技術は、この技術分野の当業者に周知であるため、ここで詳細に説明する必要はない。

【0019】【システム構成】上記のように、図1は、機械学習によるビデオサマリ要約作成システム・方法の一実施例の動作を示す概略流れ図である。システムは、ビデオ入力の画像トラックおよびオーディオトラックを検査する。さらに、システムは、入力ビデオに関連する字幕があればそれも検査することが可能である。ビデオサマリ作成システム・方法は、これらの3つの入力コンポーネント、すなわち、字幕、オーディオトラック、および画像トラックの間の整理を実行することが可能である。各入力コンポーネントに対する特徴抽出および特殊な操作も実行可能である。抽出された特徴および各コンポーネント操作の出力はその後、オーディオビジュアル統合型サマリ、または、オーディオ中心型サマリもしくは画像中心型サマリのいずれかを作成するために、機械学習によるサマリ作成モジュールに入力される。以下の操作が一般に、入力コンポーネントのそれぞれに関して実行される。

【0020】サウンド（音）の検出と分類：音楽、拍手、叫び声、爆発、雷鳴、銃声などのような非音声サウンドからなるオーディオセグメントを検出する。それらを、それぞれがコヒーレントな音響プロファイルを有するサウンドユニットに分類化する。これらのユニットを、それらの内容に従って分類する。各サウンドユニットに対して、以下のデータ、すなわち、オーディオトラック内でのそのサウンドユニットの開始時刻コード、そのサウンドユニットの継続時間、およびサウンドユニットのカテゴリあるいはタイプを出力する。

【0021】音声認識：サウンド検出・分類モジュールによって検出された非音声オーディオセグメントを取り除く。残りのオーディオセグメントに対して音声認識を実行して音声トランスクリプトを生成する。それぞれの認識語、オーディオトラック内でのその開始時刻コード、およびその継続時間を出力する。

【0022】字幕と音声トランスクリプトの整理：字幕と、音声認識器からの音声トランスクリプトとの間の整理を実行する。字幕は、タイピングミスを含むことがあり、音声認識器からの音声トランスクリプトは認識エラーを含むことがある。字幕と音声トランスクリプトの間の整理は、音声トランスクリプトの精度を改善するために有効である。

【0023】音声ユニットと特徴ベクトル生成：整理した音声トランスクリプトに基づいて音声操作ユニット

を生成し、各音声ユニットに対して特徴ベクトルを生成する。例えば、意味的な音声ユニットとしては、1個の単語、1個の句、1個の文、あるいはその他の意味的な音声内容を有するセグメントがある。

【0024】音声ユニット重要度ランク付け：各音声ユニットの重要度ランクを計算する。この重要度ランク付けは、例えば、米国特許仮出願第60/254, 535号（出願日：2000年12月12日）、発明の名称：“Text Summarization Using IR Technique And Singular Value Decomposition”）、および、米国特許出願第09/817, 591号（出願日：2001年3月26日）、発明の名称：“Generic Text Summarization Using Relevance Measure and Latent Semantic Analysis”）に記載されているような、当業者に知られた方法を利用することが可能である（本出願人による特願2001-356813号を参照）。

【0025】話者ID抽出：音声ユニット・特徴ベクトル生成モジュールから音声ユニットを受け取る。それぞれの音声ユニットに対して、話者の識別を決定する（すなわち、話者識別抽出）。

【0026】画像セグメント分割化：画像トラックを、それぞれがコヒーレントな画像プロフィールおよび動きプロフィールを有する個々の画像セグメントに分割化する。得られた画像セグメントは、画像操作ユニットとして使用可能である。

【0027】画像特徴ベクトル生成：画像特徴を抽出し、各画像セグメントに対して特徴ベクトルを形成する。特徴ベクトルを形成するためには、画像セグメント内容の何らかの側面を捕捉する任意の画像特徴が使用可能である。

【0028】顔ID抽出：それぞれの画像セグメントに人間の顔が含まれていれば、それを検出し識別する。

【0029】上記の操作が実行された後、出力は、機械学習によるサマリ作成モジュールに供給され、そこで、オーディオおよびビジュアルサマリが、前述のような機械学習フレームワークを用いて作成される。オーディオサマリ作成プロセスには、それぞれの音声あるいはサウンドユニットXに対して、そのユニットがオーディオサマリに含まれるのに十分な重要性を有する確率 $P(y|X)$ を計算することが含まれる。上記のように、それぞれの音声あるいはサウンドユニットに関連する以下の特徴が、機械学習フレームワークで使用可能である。すなわち、その特徴とは、音声ユニットの開始時刻コード、継続時間、および重要度ランク、サウンドユニットの開始時刻コード、継続時間、およびカテゴリ、ならびに、対応する画像の顔識別、および画像特徴ベクトルである。それぞれの音声あるいはサウンドユニットXに対する確率 $P(y|X)$ が計算された後、オーディオサマリがユーザ指定の長さ L_m に達するまで、確率 $P(y|X)$ の降順に音声ユニットあるいはサウンドユニットを

選択することによって、オーディオサマリが作成される。

【0030】他方、ビジュアルサマリ作成は、上記の操作で作成された画像セグメントを操作ユニットとして使用する。ビジュアルサマリ作成プロセスは、同様に、機械学習フレームワークを用いて、例えば各画像セグメントSに対して、その画像セグメントがビジュアルサマリに含まれるのに十分な重要性を有する確率 $P(y|S)$ を計算する。上記のように、例えば、各画像セグメントSに関連する以下の特徴が考えられる。すなわち、その特徴とは、長さ（すなわち、連続する、順次的な、あるいはその他の関連するフレームの個数）、画像特徴ベクトル、その画像セグメントに描画された人物あるいは顔の識別、黒フレームや画像ロゴなどのような特殊なフレームの存在、人間および物体（オブジェクト）の動き、ズームやパンなどのようなカメラの動き、対応する音声ユニットおよびサウンドユニット、ならびに、対応する音声ユニットに関連する話者の識別である。各画像セグメントSに対する確率 $P(y|S)$ が計算された後、ビジュアルサマリがユーザ指定の長さ L_m に達するまで、確率 $P(y|S)$ の降順に画像ユニットを選択することによって、ビジュアルサマリが作成される。

【0031】ビジュアルサマリは、必ずしも、それぞれの選択された画像セグメントを最初から最後まで含むことは必要でない。もとのビデオ番組を構成する画像セグメントの平均時間長は長い、ユーザ指定のサマリ長 L_m は短い場合、ビジュアルサマリはほんの2、3個の画像セグメントによって構成されることになるため、もとのビジュアル内容の大幅な喪失につながる可能性がある。ユーザがビジュアルサマリ作成結果に影響を及ぼすことを可能にするため、ユーザは、サマリ長 L_m のみならず、最小再生時間 L_{min} をも指定するように要求されることも可能である。 L_{min} は、全体のビジュアル内容を理解するためにユーザがどのくらいの長さの時間を使いたいかを示す一、 L_{min} は、幅指向ビジュアルサマリと深さ指向ビジュアルサマリ間の選択権をユーザに提供する。例えば、小さい L_{min} の値は、多数の短い画像セグメントからなる幅指向のビジュアルサマリを生成するために用いられる。他方、大きい L_{min} の値は、少数の長い画像セグメントからなる深さ指向のビジュアルサマリを生成するために用いられる。

【0032】ユーザが L_{min} および L_m を指定した後、ビジュアルサマリ内に含まれることが可能な画像セグメントの総数は、 $C = \min(L_m / L_{min}, |\Omega|)$ に等しい。ただし、 $|\Omega|$ は、もとのビデオ内の画像セグメントの総数を表す。さらに、それぞれの選択された画像セグメントに割り当てられることが可能な時間長は、 $L = L_m / C$ に等しい。この状況において、ビジュアルサマリは、確率 $P(y|S)$ の降順にC個の画像セグメントを選択し、そのC個の画像セグメントのそれぞれ

の最初の1秒間をとり、それらを時間の昇順に連結することによって作成される。

【0033】オーディオサマリとビジュアルサマリの間の整列] オーディオおよびビジュアルサマリが作成された後、解決すべき最後の問題は、どのようにしてこれらの2つのサマリを同期するかである。オーディオトラックAおよび画像トラックIからなるビデオシーケンスを $V = (I, A)$ とする。Vのオーディオサマリは、 $A_{sm} = \{A(t_i, \tau_i) \in A \mid i=1, \dots, N(A_{sm})\}$ と表される。ただし、 $A(t_i, \tau_i)$ は、時刻 t_i に開始し時間 τ_i だけ継続するオーディオセグメントを表し、 $N(A_{sm})$ は、 A_{sm} を構成するオーディオセグメントの個数を表す。 A_{sm} 内のすべてのオーディオセグメントは、それらの開始時刻 t_i の昇順に配列される。同様に、Vのビジュアルサマリは、 $I_{sm} = \{I(t_i, \tau_i) \in I \mid i=1, \dots, N(I_{sm})\}$ と表され、すべてのコンポーネントはそれらの開始時刻の昇順にソートされる。

【0034】上記のように、オーディオ中心型および画像中心型サマリは、両者の問題を最小にする。すなわち、同期は、単に、もとのビデオ番組から、画像またはオーディオのそれぞれの対応部分をとることによって実現可能である。オーディオ中心型サマリについては、 $A(t_i, \tau_i) \in A_{sm}$ の場合、 $I(t_i, \tau_i) \in I_{sm}$ である。画像中心型サマリについては、 $I(t_i, \tau_i) \in I_{sm}$ の場合、 $A(t_i, \tau_i) \in A_{sm}$ である。オーディオビジュアル統合型サマリを作成するためには、オーディオサマリとビジュアルサマリが機械学習フレームワークを用いて別個に作成されるため、それぞれのオーディオセグメント $A(t_i, \tau_i) \in A_{sm}$ に対して、対応する画像セグメント $I(t_i, \tau_i)$ は必ずしも I_{sm} に属するとは限らず、逆も同様である。したがって、画像およびオーディオの両方の内容のカバレッジを、それらのいずれをも犠牲にせずに最大化するため、オーディオサマリとビジュアルサマリの間でゆるい整列が実行される。

【0035】オーディオビジュアル統合型サマリについては、どのオーディオ内容がどの画像内容と同期しなければならぬか、およびその逆はどうかについての、システム設計者の、またはユーザの要求すなわちプレファレンスが、あらかじめ規定された整列指定として、サマリ作成システムに提供される。例えば、同期は、以下の場合に所望され、あるいは要求される。(1) ビジュアルサマリ内の画像セグメントが人物を示しており、対応するオーディオセグメント画素の人物の音声を含む場合、画像セグメントをそのオーディオ対応部分に、またはその逆に、同期することが所望される。(2) オーディオサマリ内のオーディオセグメントが爆発かななり、対応する画像セグメントが爆発を示している場合、オーディオセグメントをその画像対応部分に、またはその逆に、同期することが所望される。(3) オーディオセグ

メントが、ある有名人の名前に言及する音声を含み、その有名人の写真が、そのオーディオセグメントの小さい時間ウィンドウ内の画像セグメントに示されている場合、オーディオセグメントを、その有名人の写真を示す画像セグメントに、またはその逆に、同期することが所望される。

【0036】一実施例によれば、オーディオビジュアル統合型サマリ作成は以下のように実行される。

【0037】上記のビジュアルサマリ作成プロセスと同様に、オーディオビジュアル統合型サマリ作成は、2つのパラメータ、すなわち、ビジュアルサマリを構成する各画像セグメントに対するサマリ長 L_{sm} 、および最小再生時間 L_{min} を指定することをユーザに要求することによって開始される。ユーザが深さ指向ビジュアルサマリと幅指向ビジュアルサマリとの間の選択をすることを可能にすることは別に、パラメータ L_{min} を導入するもう1つの目的は、オーディオサマリとビジュアルサマリとの間の部分的整列を実現することである。整列の主な目標は、オーディオビジュアル統合型サマリがなめらかで自然に見えるようにし、もとのビデオのオーディオおよびビジュアルの両方の内容のカバレッジを、それらのいずれをも犠牲にすることなく、最大化することである。

【0038】例えば、ニュース番組では、ニュースアンカーやレポーターによって話される文章は、ニュース記事の重要な内容を伝えている可能性が高く、オーディオサマリに含まれる高い確率が与えられる。このような文章の対応する画像部分は、スタジオのアナウンサーや現場のレポーターの接写である。オーディオサマリ内のそれぞれの話された文が、対応する画像部分とよく整列している場合、結果は、ほとんどアナウンサーやレポーターからなる画像部分を有するビデオサマリとなる。このようにして作成されるサマリは、自然でなめらかなものに見えるかもしれないが、このような自然さおよびなめらかさは、画像内容の相当な犠牲によりもたらされたものである。完全な整列により引き起こされるこの問題を解決するため、オーディオサマリとビジュアルサマリの間で、以下の部分的整列操作が代わりに行われる。

【0039】1. 上記のように、オーディオサマリは、確率の降順で、音声またはサウンドユニットを選択することによって作成される。

【0040】2. オーディオサマリ内の各コンポーネント $A(t_i, \tau_i)$ に対して、対応する画像セグメント $I(t_i, \tau_i)$ の内容をチェックする。 $A(t_i, \tau_i)$ 、 $I(t_i, \tau_i)$ のペアが、システムに提供されたあらかじめ規定された整列要件を満たす場合、時刻 t_i から L_{sm} 秒間、 $A(t_i, \tau_i)$ を $I(t_i, \tau_i)$ と整列させる。そうでない場合は、 $A(t_i, \tau_i)$ に対して整列操作を実行しない。以下の記述において、時刻 t_i を整列点という。

【0041】3. ステップ2で整列点が識別された後、

ビデオサマリ全体は、いくつかの時間パーティションに分割される。2つの隣合う整列点 t_i , t_{i+1} に対して、期間 (t_i, t_{i+1}) に対するビジュアルサマリを作成するために、以下の操作を実行する。

【0042】 a. 期間 $(t_i, t_i + L_{min})$ を $1 (t_i, L_{min}) \in 1 (t_i, \tau_i)$ で満たす。これは、 $A(t_i, \tau_i)$ と $1 (t_i, \tau_i)$ の間の部分的整列を行う。

【0043】 b. 期間 $(t_i + L_{min}, t_{i+1})$ に対するビジュアルサマリを作成するため、この期間に入る画像セグメントの集合 Θ を求める。この期間に含まれることが可能な画像セグメントの総数は、 $C = \min((t_{i+1} - t_i - L_{min}) / L_{min}, |\Theta|)$ に等しい。ただし、 $|\Theta|$ は、集合 Θ 内の画像セグメントの個数を表す。さらに、それらの画像セグメントに割り当てられることが可能な時間長は、 $L = (t_{i+1} - t_i - L_{min}) / C$ に等しい。 Θ から、最も高い確率を有する C 個の画像セグメントを選択し、その C 個の画像セグメントのそれぞれの最初の L 秒間をとり、それらを時間の昇順に連結することによって、この期間に対するビジュアルサマリを作成する。

【0044】 [ビデオサマリ作成の代替システム・方法] 上記のように、機械学習フレームワークに基づくビデオサマリ作成のシステムおよび方法は、人間の専門家が前もって作成した十分な数のサンプルビデオサマリからなるトレーニングデータが必要とする。機械学習によるサマリ作成のシステムおよび方法は、専門家のサンプルビデオサマリから学習すること、および、サンプルビデオサマリに示される挙動を模倣することによってビデオサマリを作成することが可能である。しかし、場合によっては、専門家により作られたサンプルビデオサマリを得ることが高価すぎることや非常に困難なことがある。このような場合、トレーニングデータが必要としないシステムおよび方法を提供することが好ましい。

【0045】 図2は、トレーニングサンプルを必要としない代替的なビデオサマリ作成システム・方法の一実施例の動作を示す概略流れ図である。図2からわかるように、この代替システムは、前述の機械学習によるシステムのものと非常に類似した流れ図を有する。したがって、これから説明する代替システム・方法でも、第1実施例の場合と同様に、オーディオ中心型、画像中心型、またはオーディオビジュアル統合型のサマリを得ることが可能である。図2の流れ図において、この代替システムの、以下のモジュール以外はすべて、図1に示した対応するモジュールと同一である。

【0046】 ビジュアル内容による画像セグメントクラスタ化：第1実施例と同様の画像セグメント分節化に加えて、画像セグメントを、それらのビジュアル類似度および動的レベルに基づいてクラスタ化する。このクラスタ化は、例えば、Y. Gong and X. Liu, "Video Summarization Using Singular Value Decomposition", in Pro-

ceedings of IEEE International Conference of Computer Vision and Pattern Recognition (CVPR'00)、に記載されているものや、Y. Gong and X. Liu, "Summarizing Video By Minimizing Visual Content Redundancy", in Proceedings of IEEE International Conference of Multimedia and Expo (ICME'01)、に記載されているもののような方法を使用可能である。各フレームクラスは、同じクラス内のすべての画像セグメントが互いに視覚的に類似しているような1個以上の画像セグメントからなる。

【0047】 すべての特徴抽出操作がそれぞれのモジュールによって実行された後、出力は、オーディオビジュアルサマリ作成モジュールに供給され、そこで、オーディオサマリもしくはビジュアルサマリのいずれか、またはオーディオビジュアル統合型サマリが以下で説明するようによって作成される。

【0048】 前述のシステムにおける機械学習によるビデオサマリ作成モジュールとは異なり、この場合のオーディオビジュアルサマリ作成モジュールは、それぞれの音声またはサウンドユニット X に対する確率 $P(y|X)$ も、それぞれのフレームクラス S に対する確率 $P(y|S)$ も計算しない。代わりに、オーディオサマリがユーザ指定の長さ L_{min} に達するまで、音声ユニットを(音声ユニット重要度ランク付けモジュールから受け取った)その重要度ランクの降順に選択することによって、オーディオサマリを作成する。サウンドユニットは、例えば発見的ルールを用いて、ランク付けされ、オーディオサマリに含められるかどうか選択される。前述のように、音声ユニットの重要度ランク付けは、例えば、米国特許仮出願第60/254,535号(出願日:2000年12月12日、発明の名称:"Text Summarization Using IR Technique And Singular Value Decomposition")、および、米国特許出願第09/817,591号(出願日:2001年3月26日、発明の名称:"Generic Text Summarization Using Relevance Measure and Latent Semantic Analysis")に記載されているような、当業者に知られた方法を利用することが可能である(本出願人による特願2001-356813号を参照)。さらに、重要度ランク付けは、発見的ルールと上記の方法との組合せを用いて決定することも可能である。例えば、このような発見的ルールは、重要な人物によって話された特定の語句や、爆発、自然災害、暴行などのような特別な事件を含む重要な画像セグメントに対応する音声ユニットに、より高いランクを与える。

【0049】 ビジュアルサマリを作成するため、代替システムもまた、2個のパラメータ L_{min} 、 L_{max} のユーザによる指定を必要とする。ここでは、各フレームクラス S の重要度をランク付けするために、発見的ルールが使用可能である。一般に、ビジュアル内容サマリ作成に対するシステム設計者の、またはユーザの知識、要望、

あるいはプレファレンスを反映する任意のルールが、発見的ルールとして使用可能である。例えば、このような発見的ルールは、特定の画像特徴や、有名人や、会社ロゴなどのような特別のフレームを含むフレームクラス、人間や物体の動きや、ズーム、パンなどのようなカメラの動きを有するフレームクラス、あるいは、対応する音声ユニットが重要であるか、重要な人物によって話されているフレームクラスに、より高いランクを与える。

【0050】さらに、各フレームクラスの全時間長（構成する各画像セグメントの継続時間の和）もまた、ランク決定の過程で使用可能である。各フレームクラスは複数の画像セグメントからなることがあるため、1つのフレームクラスが選択された後、ビジュアルサマリを作成するためには、そのクラス内のどの画像セグメントを使用すべきかを決定することが依然として必要である。画像セグメント選択の助けとなる情報が他にない場合、最も直接的な選択方法は、クラス内で最長の画像セグメントを選択することとなるであろう。この理由は、同じクラス内の画像セグメントはすべて視覚的に類似しているため、最長の画像セグメントは、最も完全なものであり、最もよくクラス全体を代表するからである。この状況において、ビジュアルサマリ内に含まれることが可能な画像セグメントの総数Cと、それぞれの選択された画像セグメントに割り当てられることが可能な時間長Lは、2個のパラメータ L_{min} 、 L_{max} を利用した前述のと同じ式を用いて計算可能である。ビジュアルサマリは、C個のフレームクラスをそれらの重要度ランクの降順に選択し、そのC個のフレームクラスのそれぞれから最長の画像セグメントの最初のL秒間をとった後、それらを時間の昇順に連結することによって、作成することができる。

【0051】オーディオビジュアル統合型サマリについては、前述の機械学習によるシステムと同様に、代替システムもまた、どのオーディオ内容がどの画像内容と同期しなければならないかを示す整列指定と、パラメータ L_{min} 、 L_{max} のユーザによる入力が必要とする。オーディオサマリが作成された後、オーディオサマリ内の各コンポーネントA(t_i , τ_i)は、その画像対応部分I(t_i , τ_i)とともに検査され、A(t_i , τ_i)、I(t_i , τ_i)のペアがあらかじめ規定された整列要件を満たすかどうか調べられる。この検査は、ビデオサマリ全体をいくつかの時間パーティションに分割する整列点の集合を識別する。どのフレームクラスが、および、選択されたフレームクラス内のどの画像セグメントが、どの時間パーティションを満たすために使用されるべきかを決定しなければならない。この整列操作は、以下の2つの主なステップからなる。

【0052】1. オーディオサマリ内の各コンポーネントA(t_i , τ_i)に対して、対応する画像セグメントI

(t_i , τ_i)の内容をチェックする。A(t_i , τ_i)、I(t_i , τ_i)のペアが、あらかじめ規定された整列要件を満たす場合、時刻 t_i から L_{min} 秒間、A(t_i , τ_i)をI(t_i , τ_i)と整列させる。そうでない場合は、A(t_i , τ_i)に対して整列操作を実行しない。ここで、時刻 t_i を整列点という。

【0053】2. ステップ1ですべての整列点で識別された後、ビデオサマリ全体は、いくつかの時間パーティションに分割される。（ビジュアル内容によるフレームクラス化モジュールから得られる）クラス集合からのフレームクラスを割り当てることによって、それぞれのパーティションを満たす。この割り当ては、以下の2つの制約に適合しなければならない。

【0054】a. 単一割当て制約： 各フレームクラスは、ただ1つの時間スロット割当てを受け取ることができ。

【0055】b. 時間順序制約： ビジュアルサマリを構成するすべての画像セグメントの時間順序は維持されなければならない。

【0056】以下で、上記の整列操作のステップ2の実現法について説明する。ビデオサマリ全時間長 L_{max} が整列点によってP個のパーティションに分割され、パーティションの時間長が T_i （図3A参照）であると仮定した場合、各時間スロットは少なくとも L_{min} 秒間の長さでなければならないため、パーティション i は、

$$S_i = \lceil T_i / L_{min} \rceil$$

個の時間スロットを提供することが可能であり、したがって利用可能な時間スロットの総数は $S_{total} = \sum_{i=1}^P S_i$ となる。ここで、問題は次のようになる。ビデオサマリ全体の全部でO個のフレームクラスと S_{total} 個の時間スロットが与えられた場合に、上記の2つの制約を満たすように、フレームクラスと時間スロットの間の最適なマッチングを決定せよ。

【0057】若干の再定式化によって、今述べた問題を、最大2部マッチング問題に変換することができる。頂点の有限集合をVとし、V上の辺集合をEとする無向グラフを $G = (V, E)$ で表す。2部グラフとは、無向グラフ $G = (V, E)$ であって、 V が U 、 v 、 $U \cap v \in E$ ならば $U \cap L$ かつ $v \in R$ かつ $U \cap R$ かつ $v \in L$ のいずれかが成り立つような2つの集合LおよびRに分割可能であるようなもののことである。すなわち、すべての辺は、2つの集合LとRの間をつなぐ。マッチングとは、辺の部分集合 $M \in E$ であって、 $u \in L$ かつ $v \in R$ である任意の頂点対 (u, v) に対して、 M の高々1つの辺が u と v の間を連結するようなもののことである。

【0058】最大マッチングとは、マッチング M であって、任意のマッチング M' に対して、 $|M| \geq |M'|$ となるようなもののことである。この問題に最大2部マ

マッチングを適用するため、各頂点 $u \in U$ を用いてフレームクラスタを表し、各頂点 $v \in R$ を用いて時間スロットを表す。辺 (u, v) が存在するのは、フレームクラスタ u が、時間順序制約に違反せずに時間スロット v をとることができる場合である。フレームクラスタが、もとのビデオの前半からのものと、もとのビデオの後半からのものという複数の画像セグメントからなる場合、このフレームクラスタは、それから出た R 内の相異なる頂点に至る複数の辺を有することになる。

【0059】最大2部マッチング解は、すべてのフレームクラスタと時間スロットの間の最適割当てである。なお、最適割当ては必ずしも一意的であるとは限らない。

【0060】図3Aに、オーディオサマリとビジュアルサマリの間の整列プロセスを示す。この図において、もとのビデオ番組は70秒間の長さであり、その画像トラックは、それぞれ10秒間継続する7個の画像セグメントからなり、オーディオトラックは、それぞれ長さ10秒間の7個の語話された文からなる。ユーザは、 $L_m = 20$ 秒、および $L_w = 3$ 秒と設定している。オーディオサマリ作成は、2個の語話された文 $A(0, 10)$ および $A(30, 10)$ を選択し、ビジュアル内容によるクラスタ化は、次の2個のクラスタを生成したと仮定する：

$I(0, 10)$ からなるクラスタ1、
 $I(10, 10)$ および $I(50, 10)$ からなるクラスタ2、
 $I(30, 10)$ からなるクラスタ3、
 $I(20, 10)$ および $I(40, 10)$ からなるクラスタ4、
 $I(60, 10)$ からなるクラスタ5。

【0061】オーディオサマリが $A(0, 10)$ および $A(30, 10)$ から形成されているので、対応する画像セグメント $I(0, 10)$ および $I(30, 10)$ の内容を検査し、 $A(0, 10)$ および $A(30, 10)$ に対して整列操作が要求されるかどうかを判定する必要がある。 $I(0, 10)$ および $I(30, 10)$ は語話された文 $A(0, 10)$ 、 $A(30, 10)$ のそれぞれの話者を表示していると仮定する。その場合、整列ルールにより、 $L_m = (3)$ 秒間、 $I(0, 10)$ は $A(0, 10)$ と整列し、 $I(30, 10)$ は $A(30, 10)$ と整列することになる。 $I(0, 10)$ および $I(30, 10)$ は、一度使用されたため、これらはビジュアルサマリの他の部分で使われることはない。

【0062】これらの2つの整列法により、ビジュアルサマリの残りの期間は2つのパーティションに分割される。各パーティションは、高々2個の時間スロットを提供することが可能で7秒間継続する。整列のために3個のフレームクラスタおよび4個の時間スロットが残っているため、この整列作業に対して、図3Bに示す2部グラフがある。フレームクラスタ2は、2個の画像セグ

メント1 $(10, 10)$ および1 $(50, 10)$ からなるため、パーティション1またはパーティション2のいずれに時間スロットをとることも可能である。 $I(10, 10)$ がフレームクラスタ2から選択される場合、これはパーティション1に時間スロット2または3のいずれかをとることができる。他方、 $I(50, 10)$ が選択される場合、これはパーティション2に時間スロット5または6のいずれかをとることができる。したがって、クラスタ2から出る4本の辺、すなわち、時間スロット2への辺、時間スロット3への辺、時間スロット5への辺、および時間スロット6への辺が存在する。同様に、クラスタ4から出る4本の辺、すなわち、時間スロット2への辺、時間スロット3への辺、時間スロット5への辺、および時間スロット6への辺が存在する。

【0063】他方、フレームクラスタ5は、ただ1つの画像セグメント $I(60, 10)$ からなり、パーティション2に時間スロット5または6のいずれかをとることができる。したがって、フレームクラスタ5から出る2本の辺が存在する。

【0064】図3Bの2部グラフに対してはいくつかの可能な最大マッチング解が存在する。図4Aおよび図4Bは2つのそれぞれの解を示す。図4Aに示す解(i)では、時間スロット3が未割当てのままである。図4Bに示す解(ii)では、時間スロット5が未割当てのままである。この場合、すべてのフレームクラスタが使用されているため、複数の画像セグメントを有するフレームクラスタを用いて、空き時間スロットを満たす必要がある。解(i)(図4A)の場合、フレームクラスタ4の画像セグメント $I(20, 10)$ が、空き時間スロットを満たすために使用されなければならない。解(ii)(図4B)の場合、フレームクラスタ2の画像セグメント1 $(50, 10)$ が、空き時間スロットを満たすために使用されなければならない。

【0065】上記の例は次のことを例示している。すなわち、最大2部マッチングは、利用可能なフレームクラスタと時間スロットの間の最適なマッチングを求めるが、特に、利用可能なフレームクラスタの数より多くの利用可能な時間スロットがあるときには、一部の時間スロットを未割当てのまま残すことがある。これらの未割当て時間スロットを満たすために、単一割当て制約をゆるめ、複数の画像セグメントを有するフレームクラスタを検査し、まだ使用されていない適当なセグメントを選択することが可能である。このようにして、時間順序制約は満たされる。ゆるめられた単一割当て制約に対するそれぞれの解を図5Aおよび図5Bに示す。

【0066】なお、最大2部マッチング操作は、不正な解を生成することがある。図6Aおよび図6Bは、これらの2つの例を示す。例(i)(図6A)では、画像セグメント $I(60, 10)$ が画像セグメント $I(50, 10)$ の前に置かれているため、時間順序制約に違反して

いる。例(i i)(図6B)では、割当てはいずれの制約にも違反していないが、1(20, 10)を時間スロット2に割り当てることが、時間スロット3の割当てを不可能にしている。しかし、これらの不正な解は、これら2つの制約に照らして検査することによって容易に検出され、各パーティションにおいて時間スロットに割り当てられる画像セグメントを並べ替えることによって補正することができる。例(i i)(図6A)の場合、問題は、パーティション2に割り当てられた2個の画像セグメントを時間の昇順にソートすることによって補正することができる。例(i i)(図6B)の場合、まず、フレームクラスタ2からの画像セグメントI(10, 10)(これは、パーティション1に割り当てられることが可能な唯一の残りの画像セグメントである)を用いて空き時間スロットを満たした後に、そのパーティション内の2個の画像セグメントをソートすることによって、最終的な解に達することができる。

【0067】まとめると、整列操作のステップ2は、次のように記述することができる。

【0068】1. 整列点が識別された後、割当てのために残っているフレームクラスタおよび時間スロットの個数を決定し、それに応じて2部グラフを作る。

【0069】2. 最大2部マッチングアルゴリズムを用いて可能な解を求める。

【0070】3. 解を2つの制約について検査し、各パーティション内の画像セグメントをソートすることによって違反を補正する。

【0071】4. 未割当て時間スロットが存在する場合、単一割当て制約をゆるめ、複数の画像セグメントを有するフレームクラスタを検査し、まだ使用されていないセグメントで時間順序制約を満たす適当なセグメント*

*を選択する。

【0072】以上、好ましい実施例を参照して、本発明について詳細に説明したが、本発明の技術的範囲および技術思想の範囲内のさまざまな変形は、この技術分野の当業者には明らかである。したがって、本発明は、特許請求の範囲の技術的範囲によってのみ限定されたとみなされるべきである。

【0073】

【発明の効果】以上詳細に説明したように、本発明によれば、オーディオと画像の内容の厳密な同期が要求されないときには、オーディオビジュアル統合型サマリ作成技術を用い、オーディオ内容と画像内容の同期を要求するビデオ番組の場合には、オーディオ中心型または画像中心型のいずれかの方法を用いてサマリが作成される。これにより、入力ビデオから抽出された画像、オーディオ、およびテキスト特徴をシームレスに統合し、オーディオ中心型、画像中心型、およびオーディオビジュアル統合型の高品質のサマリを作成することができる。

【図面の簡単な説明】

【図1】本発明のオーディオビジュアルサマリ作成システム・方法の一実施例の動作を示す流れ図である。

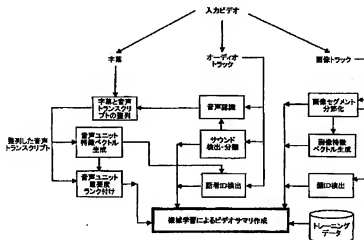
【図2】本発明のオーディオビジュアルサマリ作成システム・方法の代替実施例の動作を説明する流れ図である。

【図3】Aは、オーディオサマリとビジュアルサマリの間の整列プロセスを示す図である。Bは、その整列のためのフレームワークを示す図である。

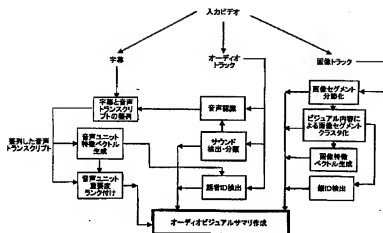
【図4】時間順序制約を満たす代替解を示す図である。

【図5】時間順序制約を満たす代替解を示す図である。
【図6】本発明の方法から得られる不正な解を示す図である。

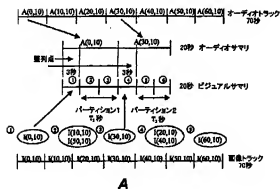
【図1】



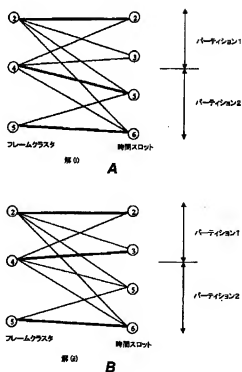
【図2】



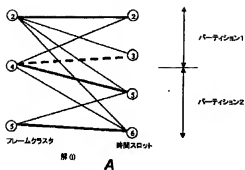
【図3】



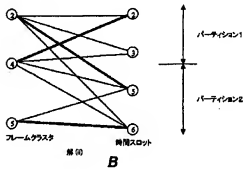
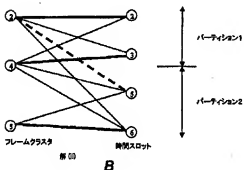
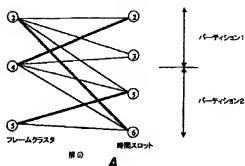
【図4】



【図5】



【図6】



フロントページの続き

(51) Int. Cl.⁷

H 0 4 N 5/91

識別記号

F I

H 0 4 N 5/91

キーワード¹ (参考)

N

C

(72) 発明者 シン リュウ

アメリカ合衆国、ニュージャージー

08540 プリンストン、4 インディペン

デンス ウエイ、エヌ・イー・シー・ユ

ー・エス・エー インク内

Fターム (参考) 5C053 FA14 GA16 GB05 JA01

5D015 AA03 AA06 FF00 GG00 HH00

KK02 LL11